

Mining Billion-Scale Graphs: Patterns and Algorithms

Christos Faloutsos and U Kang

CMU

Part 2: Algorithms

Complementary to tutorial: *Mining Billion-Scale Graphs:
Systems and Implementations*: Haixun Wang et al

Part 2: Algorithms

Outline

- Problem#1: Patterns in graphs
- Problem#2: Tools
- ➔ • Problem#3: Scalability - PEGASUS
 - Structure Analysis
 - Eigensolver
 - Graph Layout and Compression
- Conclusions

Our goal:

Open source system for mining huge graphs:

PEGASUS project (PEta GrAph mining System)

- www.cs.cmu.edu/~pegasus
- code and papers



Scalability Challenge

- The sizes of graphs are growing!

facebook

0.5 billion users
60 TBytes/day
15 PBytes/total
[Thusoo+ '10]

YAHOO!

1.4 billion web pages
6.6 billion edges
[Broder+ '04]

bing

ClickStream Data
0.26 PBytes
1 billion query-URL
[Liu+ '09]

Google

20 PBytes/day
[Dean+ '08]

Scalability Challenge

- The sizes of graphs are growing!

facebook

bing

Q: How can we handle large graphs which don't fit into the memory, or disks of a single machine?

YAHOO!

Google

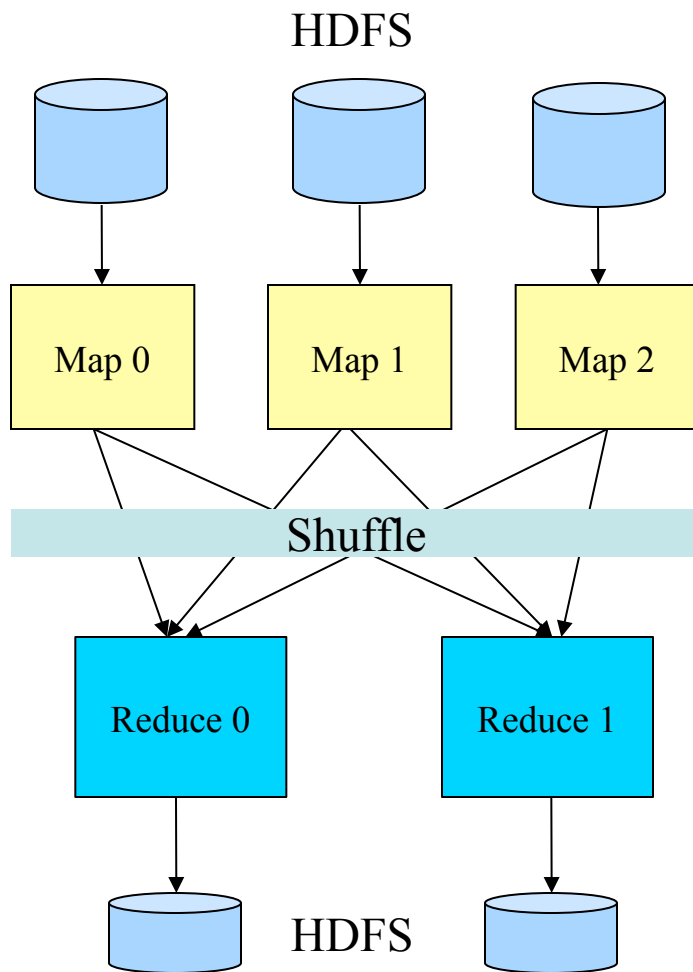
A: Parallelism, with MapReduce!

[Broder+04]

[Dean+08]

Background: MapReduce

MapReduce/Hadoop Framework



HDFS: fault tolerant, scalable, distributed storage system

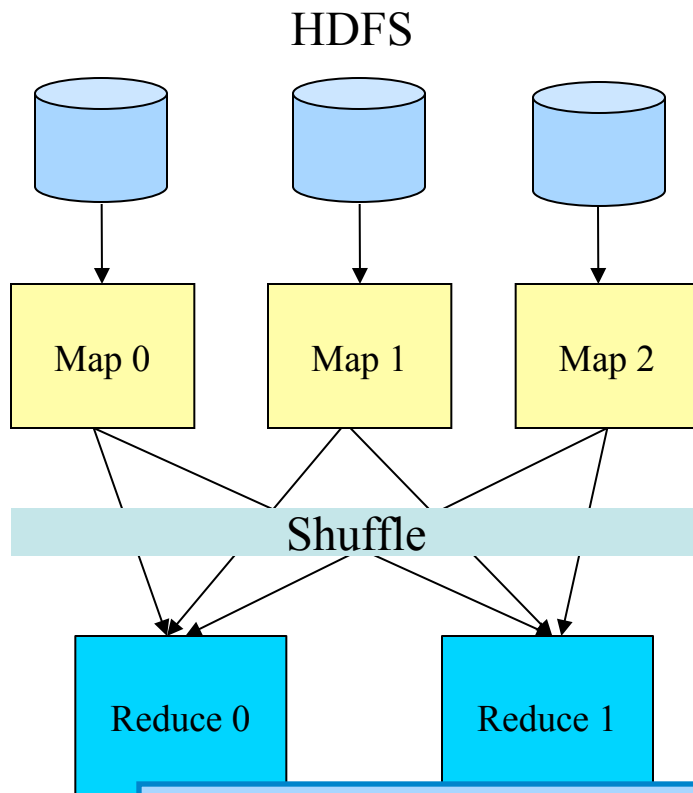
Mapper: read data from HDFS, output (k,v) pair

Output sorted by the key

Reducer: read output from mappers, output a new (k,v) pair to HDFS

Background: MapReduce♪

■ MapReduce/Hadoop Framework



HDFS: fault tolerant, scalable, distributed storage system

Mapper: read data from HDFS, output (k,v) pair

Output sorted by the key

Reducer: read output from mappers, output a new (k,v) pair to H

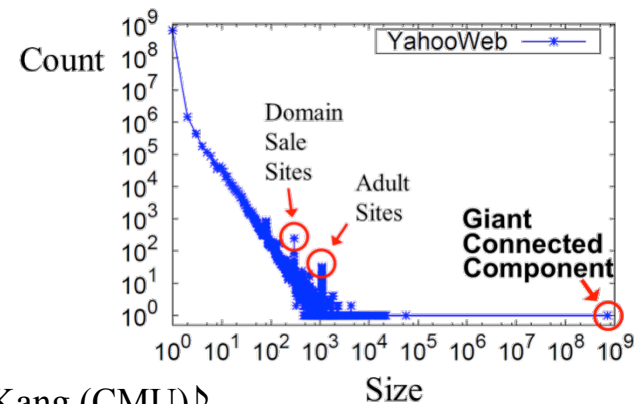
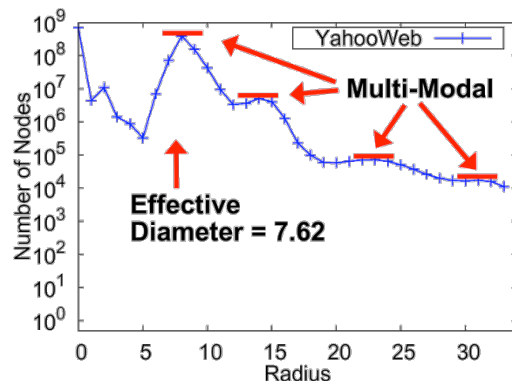
Programmers need to provide only
map() and reduce() functions

Outline

- Problem#1: Patterns in graphs
- Problem#2: Tools
- Problem#3: Scalability - PEGASUS
 - ➔ – Structure Analysis
 - Eigensolver
 - Graph Layout and Compression
- Conclusions

Structure Analysis

- How to scale-up structure analysis algorithm?
 - Q1: How to **unify** many structure analysis algorithms (connected components, PageRank, diameter/radius)?
 - Q2: How to design a **scalable** algorithm for the structure analysis?



Q1: Unifying Algorithms

- Given a graph, can we compute
 - connected components,
 - PageRank,
 - Random Walk with Restart,
 - diameter/radiuswith *one algorithm*?

Q1: Unifying Algorithms

- Given a graph, can we compute
 - connected components,
 - PageRank,
 - Random Walk with Restart,
 - diameter/radiuswith *one algorithm*?

Yes!

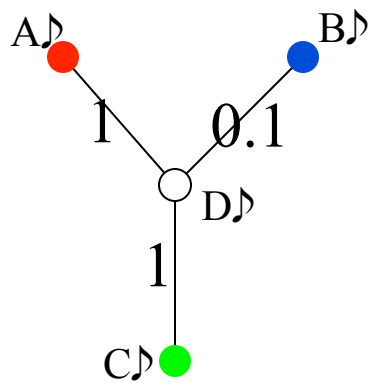
How ?

Main Idea

- GIM-V
 - Generalized Iterative Matrix-Vector Multiplication
 - Extension of plain matrix-vector multiplication
 - includes
 - Connected Components
 - PageRank
 - RWR (Random Walk With Restart)
 - Diameter Estimation

Main Idea: Intuition

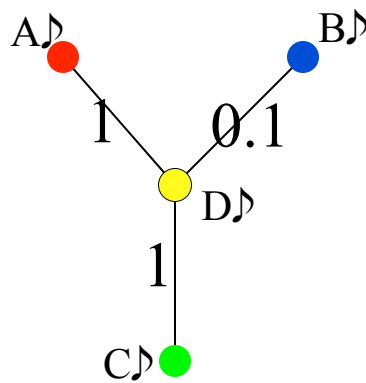
■ Plain M-V multiplication



- Weighted Combination of Colors
- ~ Message Passing

Main Idea: Intuition

■ Plain M-V multiplication



- Weighted Combination of Colors
- ~ Message Passing

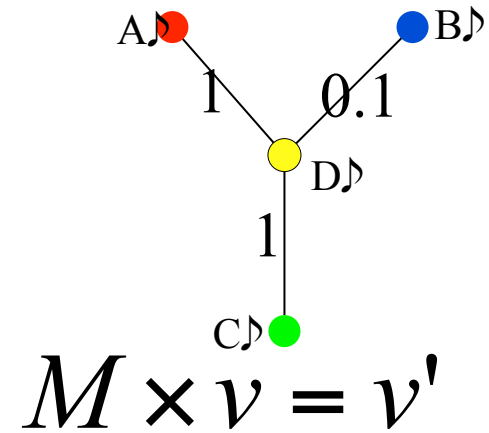
$$\begin{array}{c}
 M \\
 \begin{array}{c|c|c|c}
 & A & B & C & D \\
 \hline
 A & & & & 1 \\
 B & & & & 1 \\
 C & & & & 0.1 \\
 D & 1 & 1 & 0.1 &
 \end{array}
 \end{array}
 \times
 \begin{array}{c}
 v \\
 \begin{array}{|c|}
 \hline
 \bullet \\
 \hline
 \bullet \\
 \hline
 \bullet \\
 \hline
 \circ \\
 \hline
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 v' \\
 \begin{array}{|c|}
 \hline
 \\
 \hline
 \\
 \hline
 \\
 \hline
 \bullet \\
 \hline
 \end{array}
 \end{array}$$

$$v_4' = \sum_{i=1}^4 m_{4i} v_i$$

Main Idea: Intuition

■ Plain M-V multiplication

$$\begin{array}{c}
 M \\
 \begin{array}{c} A \rightsquigarrow B \rightsquigarrow C \rightsquigarrow D \rightsquigarrow \\
 \begin{array}{cccc}
 A \rightsquigarrow & & & 1 \\
 B \rightsquigarrow & & & 1 \\
 C \rightsquigarrow & & & 0.1 \\
 D \rightsquigarrow & 1 & 1 & 0.1
 \end{array}
 \end{array}
 \end{array}
 \times
 \begin{array}{c}
 v \\
 \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \circ \end{array}
 \end{array}
 =
 \begin{array}{c}
 v' \\
 \begin{array}{c} \\ \\ \\ \bullet \end{array}
 \end{array}$$

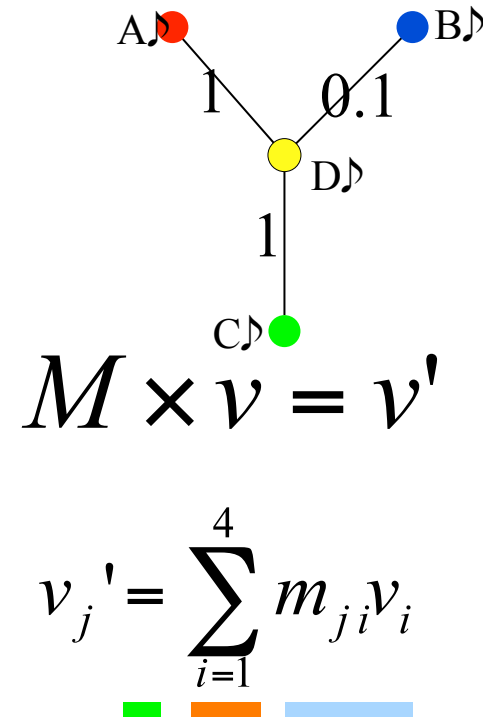


$$v_j' = \sum_{i=1}^4 m_{ji} v_i$$

Main Idea: Intuition

Plain M-V multiplication

$$\begin{array}{c}
 M \\
 \begin{array}{cccc}
 A \rightarrow & B \rightarrow & C \rightarrow & D \rightarrow \\
 A \rightarrow & & & 1 \\
 B \rightarrow & & & 1 \\
 C \rightarrow & & & 0.1 \\
 D \rightarrow & 1 & 1 & 0.1
 \end{array}
 \end{array}
 \times
 \begin{array}{c}
 v \\
 \begin{array}{c}
 \bullet \\
 \bullet \\
 \bullet \\
 \circ
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 v' \\
 \begin{array}{c}
 \\
 \\
 \\
 \bullet
 \end{array}
 \end{array}$$



Three Implicit Operations here:

multiply m_{ji} and v_i

combine2

Message sending

sum n multiplication results

combineAll


Message combination

update v_j'

assign

Main Idea

- GIM-V
 - Generalizing the three operations leads to many algorithms

operations	Standard MV	Con. Cmpt.	PageRank	RWR	Diameter
combine2	Multiply				
combineAll	Sum				
assign	Assign				

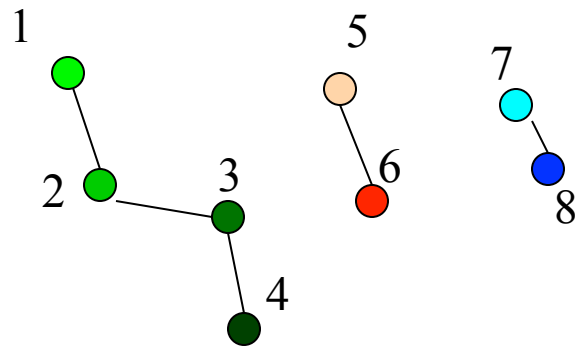
Main Idea

■ GIM-V for Connected Components

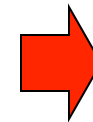
operations	Standard MV	Con. Cmpt.	PageRank	RWR	Diameter
combine2	Multiply	Bool. X			
combineAll	Sum	MIN			
assign	Assign	MIN			

Main Idea

- GIM-V for Connected Components
 - How many connected components?
 - Which node belong to which component?



Input Graph

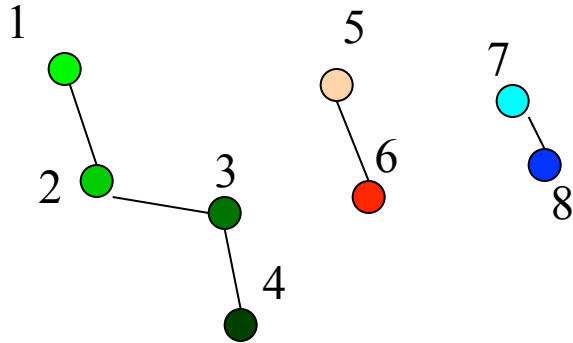


node id	component id
1	1
2	1
3	1
4	1
5	5
6	5
7	7
8	7

Output

Main Idea

■ GIM-V for Connected Components

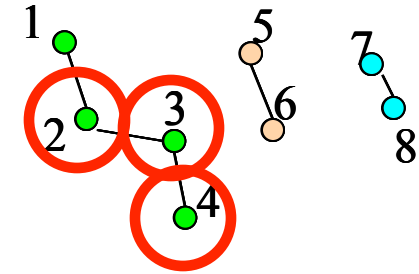


	1	2	3	4	5	6	7	8
1		1						
2	1		1					
3		1		1				
4			1					
5						1		
6					1			
7								1
8							1	

$$\begin{array}{c}
 \times G \\
 \begin{array}{c}
 \boxed{1} \\
 \boxed{2} \\
 \boxed{3} \\
 \boxed{4} \\
 \boxed{5} \\
 \boxed{6} \\
 \boxed{7} \\
 \boxed{8}
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 \boxed{1} \\
 \boxed{1} \\
 \boxed{1} \\
 \boxed{1} \\
 \boxed{5} \\
 \boxed{5} \\
 \boxed{7} \\
 \boxed{7}
 \end{array}$$

initial vector final vector

Main Idea



■ GIM-V for Connected Components♪

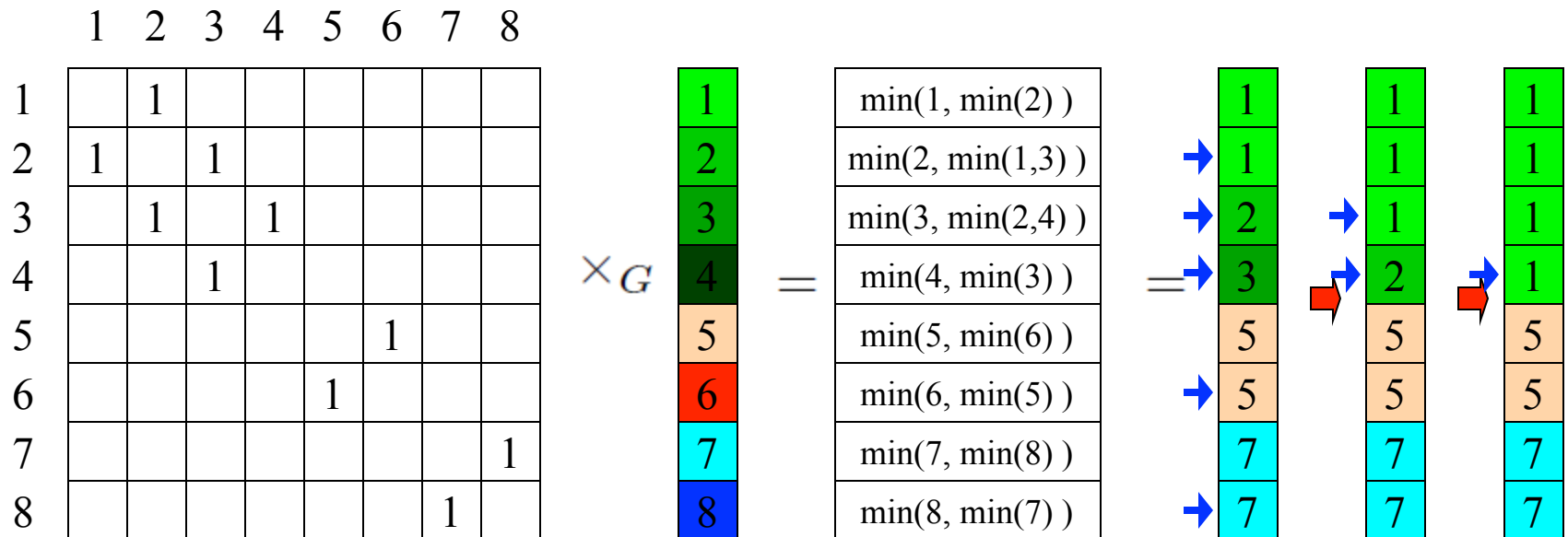
$$\text{combine2}(m_{i,j}, v_j) = m_{i,j} \times v_j$$

$$\text{combineAll}(x_1, \dots, x_n) = \min\{x_i \mid i = 1..n\}$$

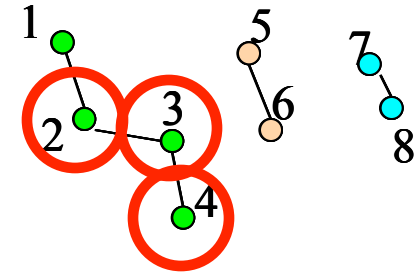
$$\text{assign}(v_i, v_{\text{new}}) = \min(v_i, v_{\text{new}})$$

“Sending Invitations”

“Accept the Smallest”



Main Idea



■ GIM-V for Connected Components♪

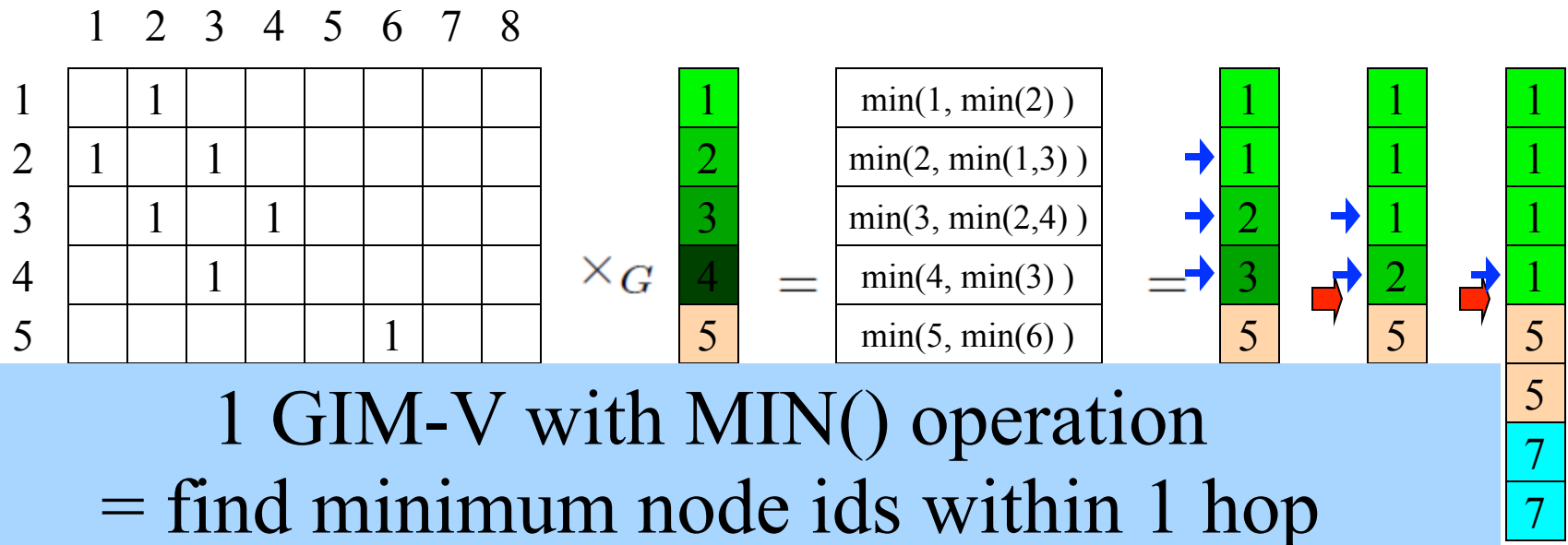
$$\text{combine2}(m_{i,j}, v_j) = m_{i,j} \times v_j$$

$$\text{combineAll}(x_1, \dots, x_n) = \min\{x_i \mid i = 1..n\}$$

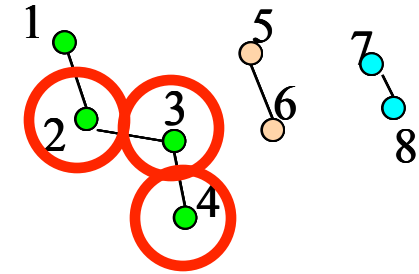
$$\text{assign}(v_i, v_{\text{new}}) = \min(v_i, v_{\text{new}})$$

“Sending Invitations”

“Accept the Smallest”



Main Idea



■ GIM-V for Connected Components♪

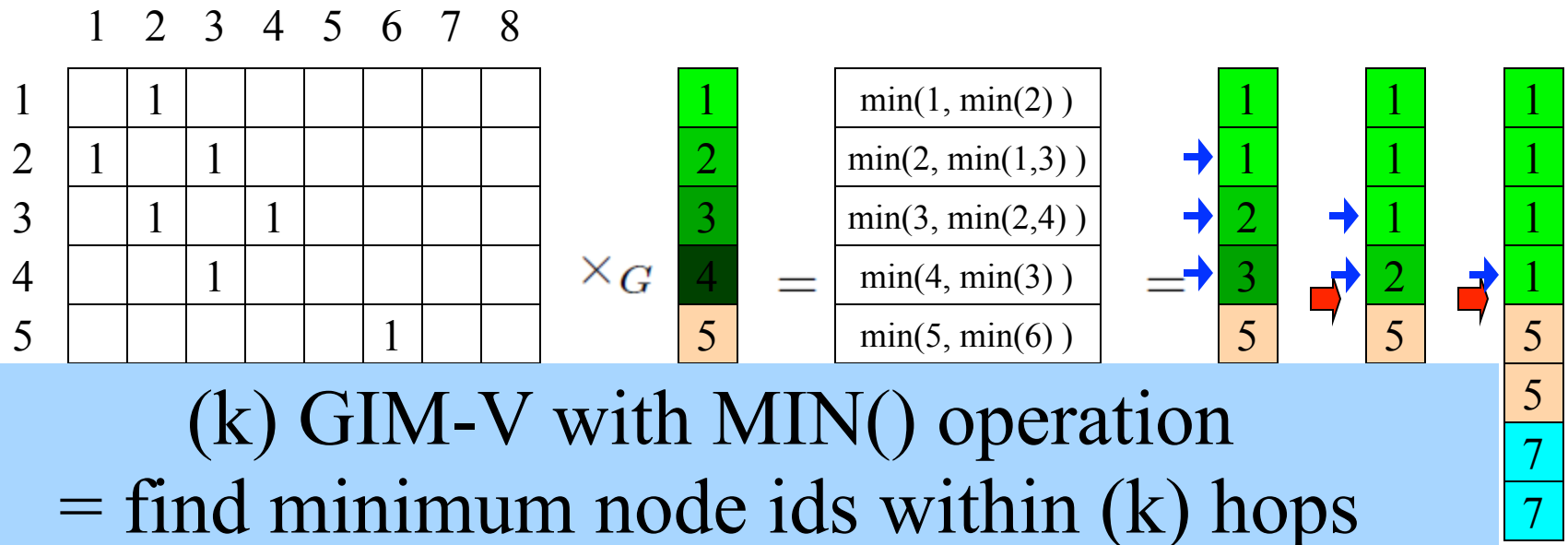
$$\text{combine2}(m_{i,j}, v_j) = m_{i,j} \times v_j$$

$$\text{combineAll}(x_1, \dots, x_n) = \min\{x_i \mid i = 1..n\}$$

$$\text{assign}(v_i, v_{\text{new}}) = \min(v_i, v_{\text{new}})$$

“Sending Invitations”

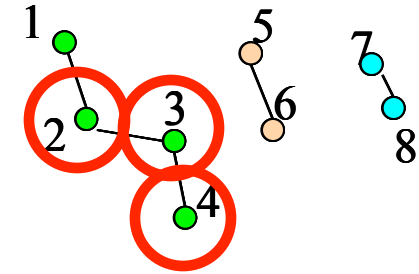
“Accept the Smallest”



(k) GIM-V with MIN() operation

= find minimum node ids within (k) hops

Main Idea



■ GIM-V for Connected Components♪

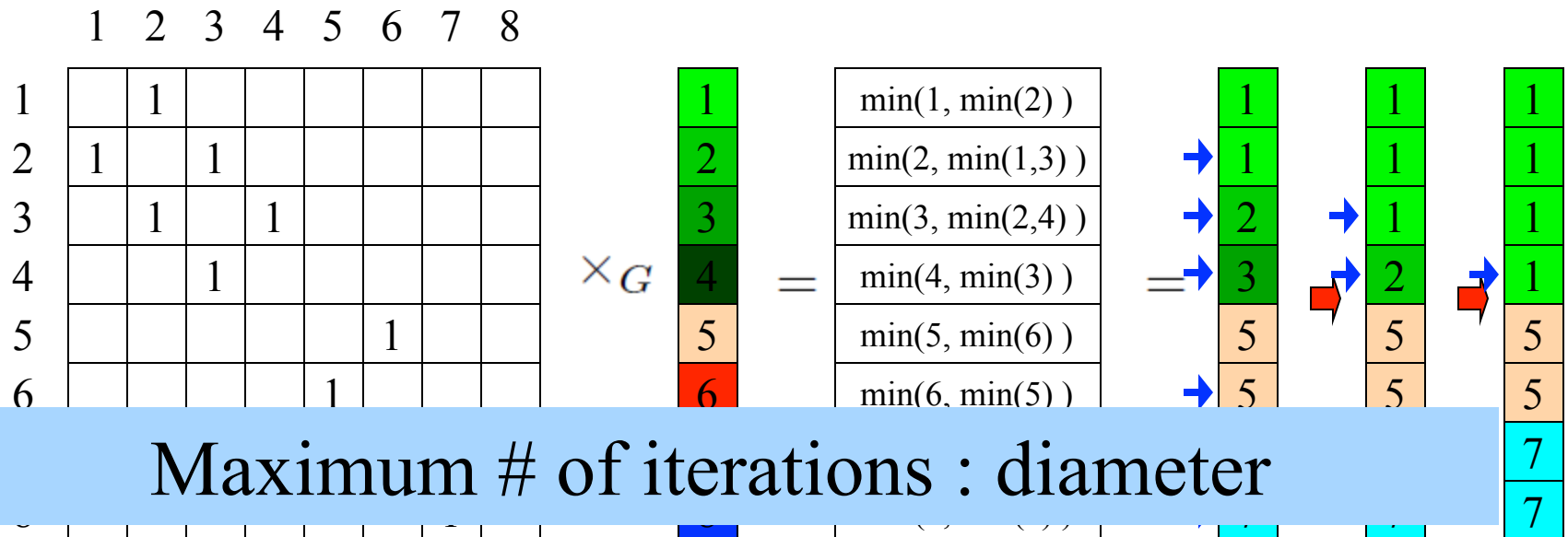
$$combine2(m_{i,j}, v_j) = m_{i,j} \times v_j$$

$$combineAll(x_1, \dots, x_n) = \min\{x_i \mid i = 1..n\}$$

$$assign(v_i, v_{new}) = \min(v_i, v_{new})$$

“Sending Invitations”

“Accept the Smallest”



Maximum # of iterations : diameter

Main Idea

■ GIM-V for PageRank

Operations	Standard MV	Con. Cmpt.	PageRank	RWR	Diameter
combine2	Multiply	Multiply	Multiply with c		
combineAll	Sum	MIN	Sum with rj prob		
assign	Assign	MIN	Assign		

Main Idea

Details

- GIM-V: PageRank
 - PageRank vector p : eigenvector of A :

$$1p = A \times p$$

$$\begin{array}{ccc} \boxed{} & \boxed{} & \boxed{} \\ \text{nx1} & \text{nxn} & \text{nx1} \end{array}$$

- Where

$$A = cE^T + (1 - c)U$$

Adjacency Matrix
All elements set to 1/n

Main Idea

Details

- GIM-V: PageRank
 - Algorithm: Power method
(many multiplications)

$$\begin{array}{ccc}
 p_{next} & \leftarrow & A \times p_{cur} \\
 \begin{array}{c} \square \\ nx1 \end{array} & & \begin{array}{cc} \square & \square \\ nxn & nx1 \end{array}
 \end{array}$$

Main Idea

Details

■ GIM-V: PageRank

$$\begin{array}{ccc}
 p_{next} & \leftarrow & E^T \times_G p_{cur} \\
 \begin{array}{c} \square \\ nx1 \end{array} & & \begin{array}{c} \square \\ nxn \end{array} \quad \begin{array}{c} \square \\ nx1 \end{array}
 \end{array}$$

Algorithm: Power method
(many multiplications)

$$combine2(m_{i,j}, v_j) = c \times m_{i,j} \times v_j$$

$$combineAll(x_1, \dots, x_n) = \frac{1-c}{n} + \sum x_i$$

$$assign(v_i, v_{new}) = v_{new}$$

Main Idea

■ GIM-V

Operations	Standard MV	Con. Cmp.	PageRank	RWR	(approx.) Diameter
combine2	Multiply	Multiply	Multiply with c	Multiply with c	Multiply bit-vector
combineAll	Sum	MIN	Sum with rj prob.	Sum with rest prob	BIT-OR()
assign	Assign	MIN	Assign	Assign	BIT-OR()

Two Restrictions on HDFS

- [R1] HDFS is location transparent
 - Users don't know which file is located in which machine
- [R2] A line is never split
 - A large file is split into pieces of a size (e.g. 256 MB)
 - Users don't know the point of the split

Q2: Fast Algorithms for GIM-V

- Given the two restrictions R1 and R2, how can we make faster algorithms for GIM-V in Hadoop?
 - Three main ideas:
 - I1) Block Multiplication
 - I2) Clustering
 - I3) Compression

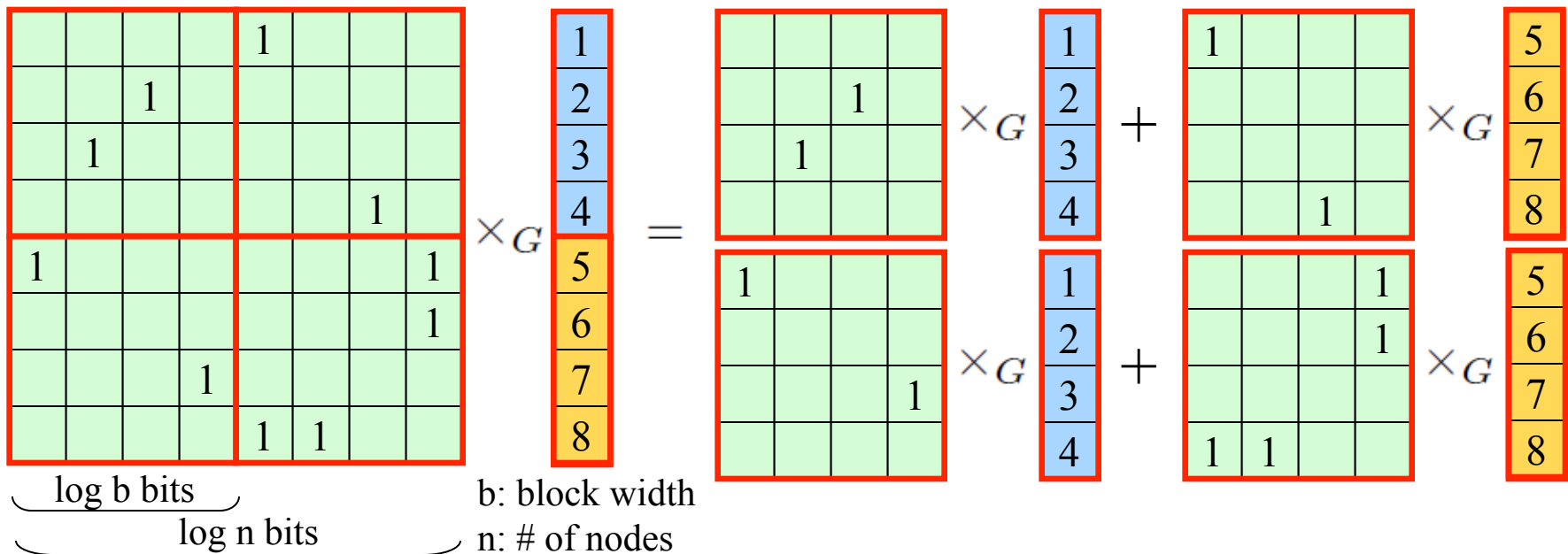
Main Idea

■ I1) Block-Method

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline & & & & 1 & & & \\ \hline & & & 1 & & & & \\ \hline & & 1 & & & & & \\ \hline & & & & & & & 1 \\ \hline 1 & & & & & & & 1 \\ \hline & & & & & & & 1 \\ \hline & & & & 1 & & & \\ \hline & & & & 1 & 1 & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline 4 \\ \hline 5 \\ \hline 6 \\ \hline 7 \\ \hline 8 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & 1 \\ \hline & & 1 & \\ \hline & & & \\ \hline 1 & & & \\ \hline & & & \\ \hline & & & 1 \\ \hline & & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline 4 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline 1 & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & 1 \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline 1 & 1 & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 5 \\ \hline 6 \\ \hline 7 \\ \hline 8 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline & & & 1 \\ \hline & & & 1 \\ \hline & & & \\ \hline & & & \\ \hline 1 & 1 & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 5 \\ \hline 6 \\ \hline 7 \\ \hline 8 \\ \hline \end{array}$$

Main Idea

■ I1) Block-Method



1. Group elements together into 1 line
2. Storage for an element: $2\log n$ bits \rightarrow $2\log b$ bits
3. Adjust the MapReduce code(block multiplication)

Main Idea

■ I1) Block-Method

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline & & & & 1 & & & \\ \hline & & & 1 & & & & \\ \hline & & 1 & & & & & \\ \hline & & & & & & & 1 \\ \hline & & & & & & 1 & \\ \hline 1 & & & & & & & 1 \\ \hline & & & & & & & 1 \\ \hline & & & & 1 & & & \\ \hline & & & & & 1 & 1 & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline 4 \\ \hline 5 \\ \hline 6 \\ \hline 7 \\ \hline 8 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & 1 \\ \hline & & 1 & \\ \hline & & & \\ \hline & & & \\ \hline 1 & & & \\ \hline & & & \\ \hline & & & 1 \\ \hline & & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline 4 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline 1 & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & 1 \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline 1 & 1 & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 5 \\ \hline 6 \\ \hline 7 \\ \hline 8 \\ \hline \end{array} + \begin{array}{|c|c|c|c|} \hline & & & 1 \\ \hline & & & 1 \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline 1 & 1 & & \\ \hline \end{array} \times_G \begin{array}{|c|} \hline 5 \\ \hline 6 \\ \hline 7 \\ \hline 8 \\ \hline \end{array}$$

Thanks to the encoding,
 - file size is decreased,
 - shuffle time is decreased.

Main Idea

Q: Can we do even better?

Main Idea

■ I2) Clustering

				1			
		1					
	1						
						1	
1							1
							1
			1				
				1	1		

Preprocess

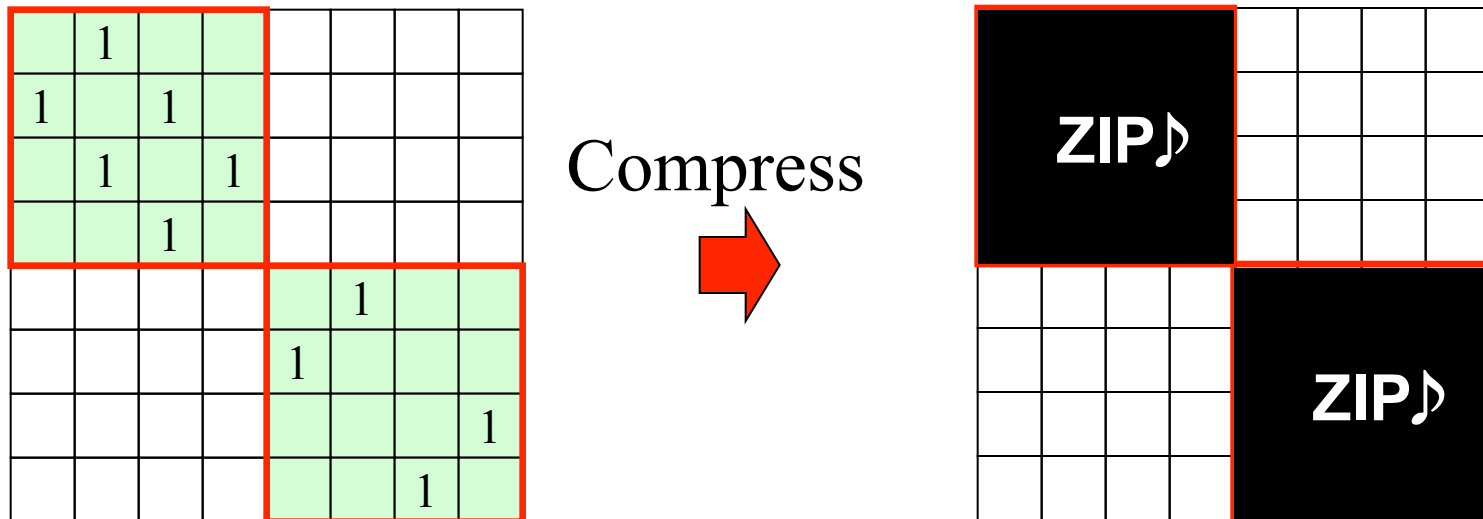


	1						
1		1					
	1		1				
		1					
					1		
					1		
							1
						1	

A: preprocessing for clustering
(only green blocks are stored in HDFS)

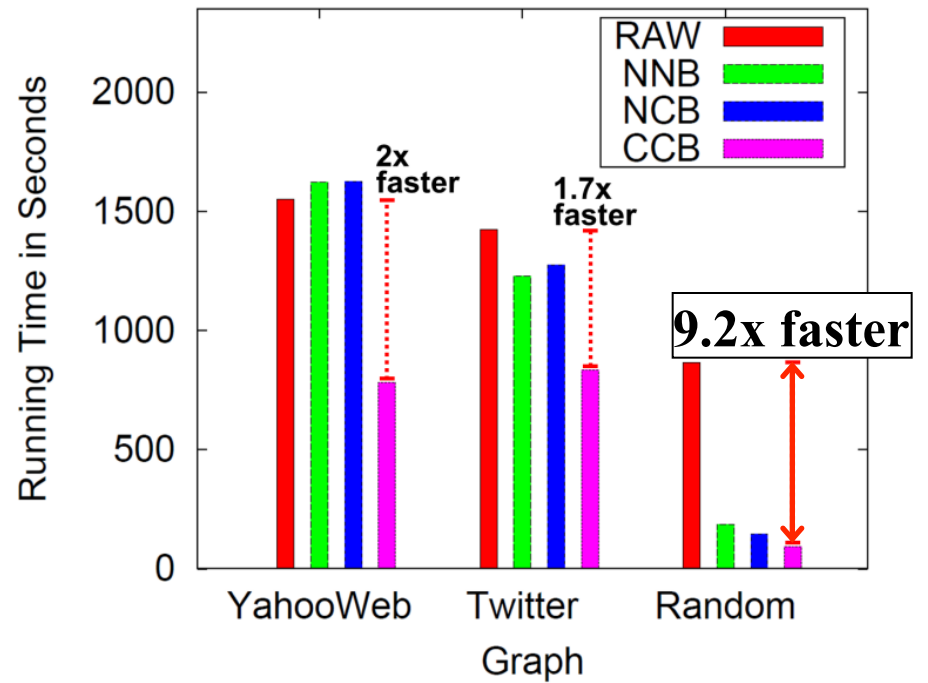
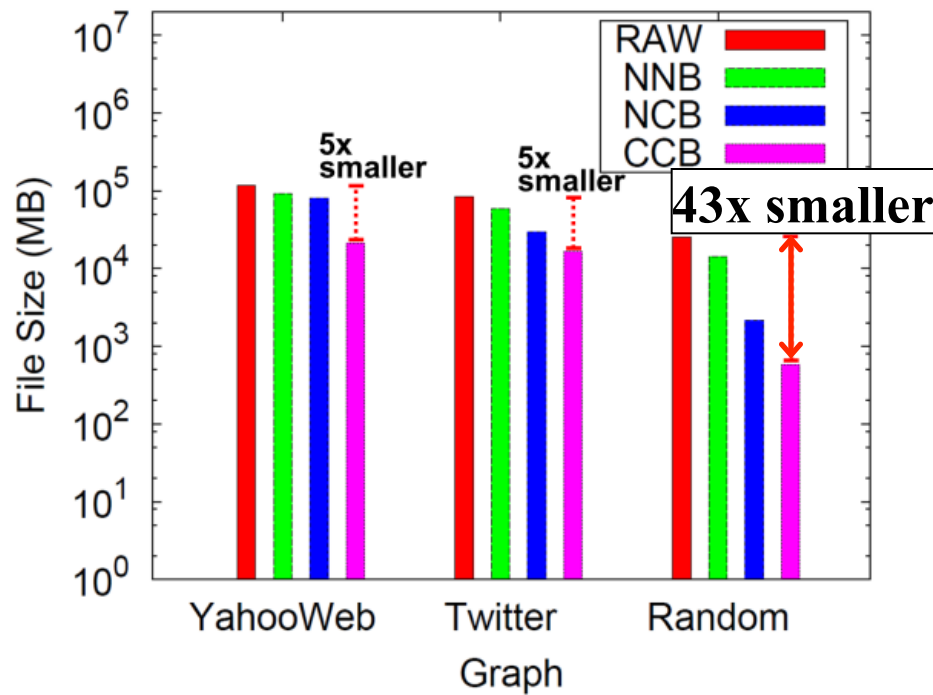
Main Idea

■ I3) Compression



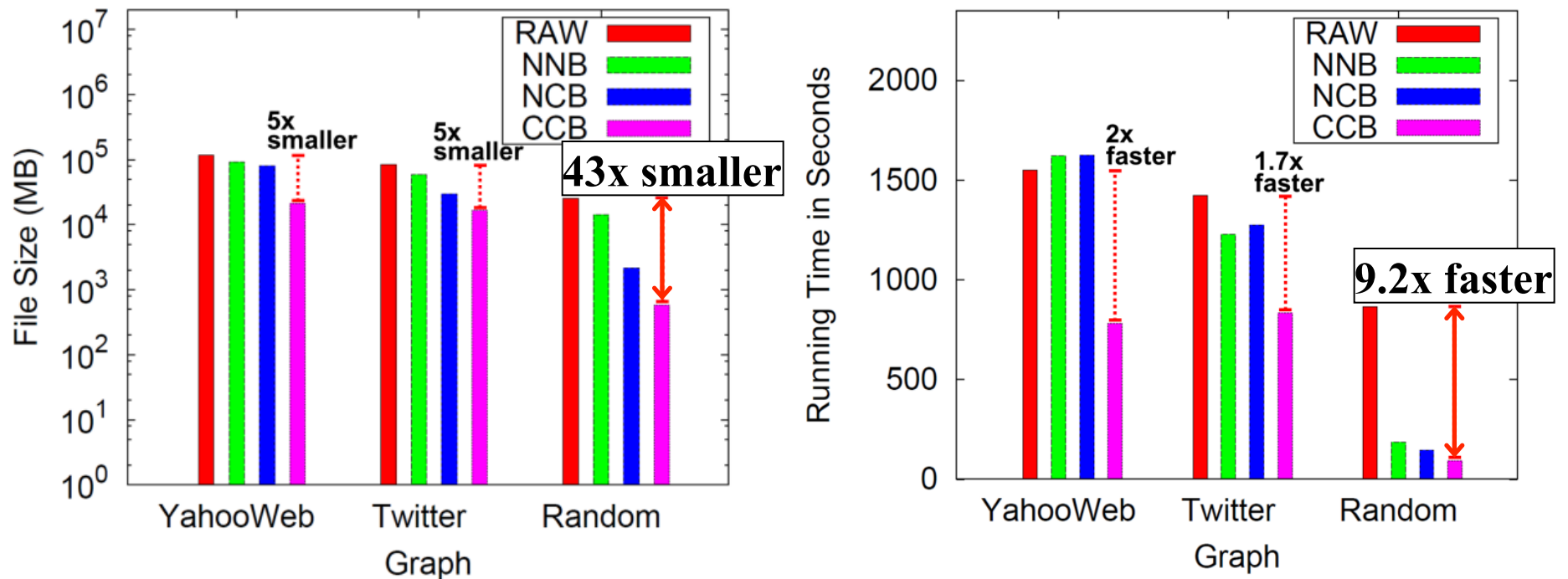
A: compress clustered blocks

Result



	Block Encoding ?♪	Compression ?♪	Clustering?♪
RAW♪	No♪	No♪	No♪
NNB♪	Yes♪	No♪	No♪
NCB♪	Yes♪	Yes♪	No♪
CCB♪	Yes♪	Yes♪	Yes♪

Result



A: Proposed Method(CCB) provides
43x smaller storage, 9.2x faster running time

Outline

- Problem#1: Patterns in graphs
- Problem#2: Tools
- Problem#3: Scalability - PEGASUS
 - Structure Analysis
 - ➔ – Eigensolver
 - Graph Layout and Compression
- Conclusions

Background: Eigensolver

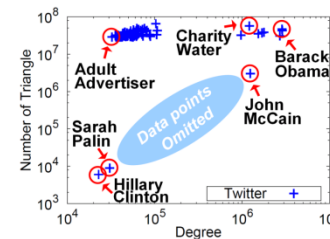
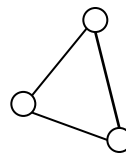
■ Eigensolver

- Given: (adjacency) matrix A ,
- Compute: top k eigenvalues and eigenvectors of A
- Application:

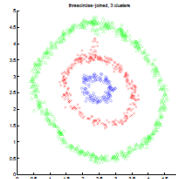
■ SVD

$$\square = \square \square \square$$

■ triangle counting



■ spectral clustering



Problem Definition

- Q4: How to design a billion-scale eigensolver?
 - Existing eigensolver: can handle millions of nodes and edges

Efficient Eigensolver

Details

■ Lanczos Iterations

$$\beta_0 = 0, q_0 = 0, b = \text{arbitrary}, q_1 = b/\|b\|$$

for $n = 1, 2, 3, \dots$

$$v = Aq_n$$

$$\alpha_n = q_n^T v$$

$$v = v - \beta_{n-1}q_{n-1} - \alpha_n q_n$$

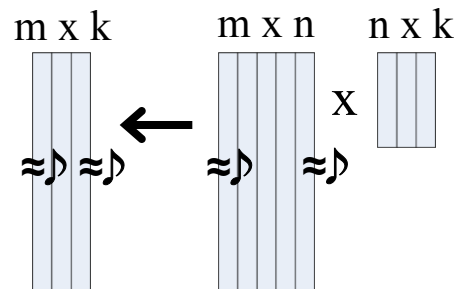
$$\beta_n = \|v\|$$

$$q_{n+1} = v/\beta_n$$

1 matrix-vector multiplication
per iteration

Proposed Method

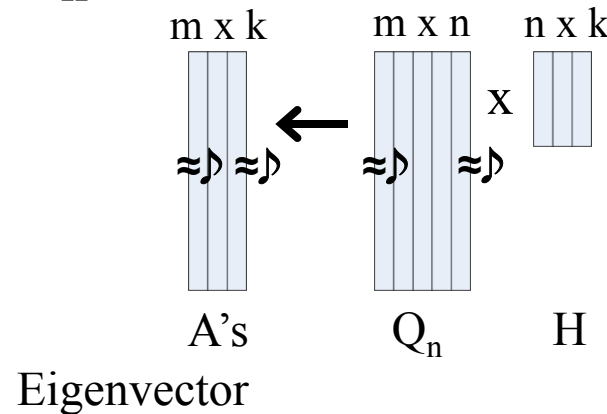
- HEigen algorithm (Hadoop Eigen-solver)
 - Selectively parallelize ‘Lanczos-SO’ algorithm
 - Block encoding
 - Exploiting skewness in matrix-matrix mult.
 - $(m \gg n > k)$



Skewed Matrix-Matrix Mult.

Details

- Multiply $Q_n^{m \times n}$ and $H^{n \times k}$ ($m \gg n > k$)



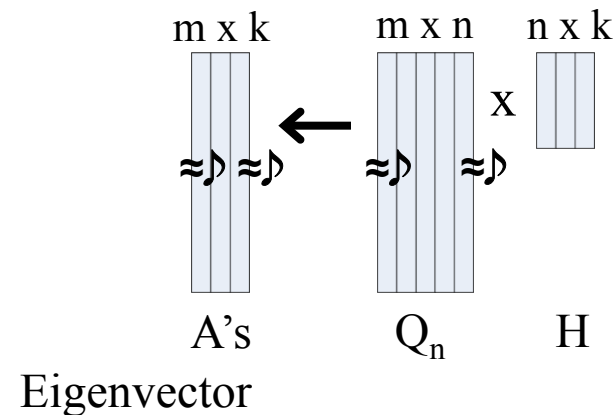
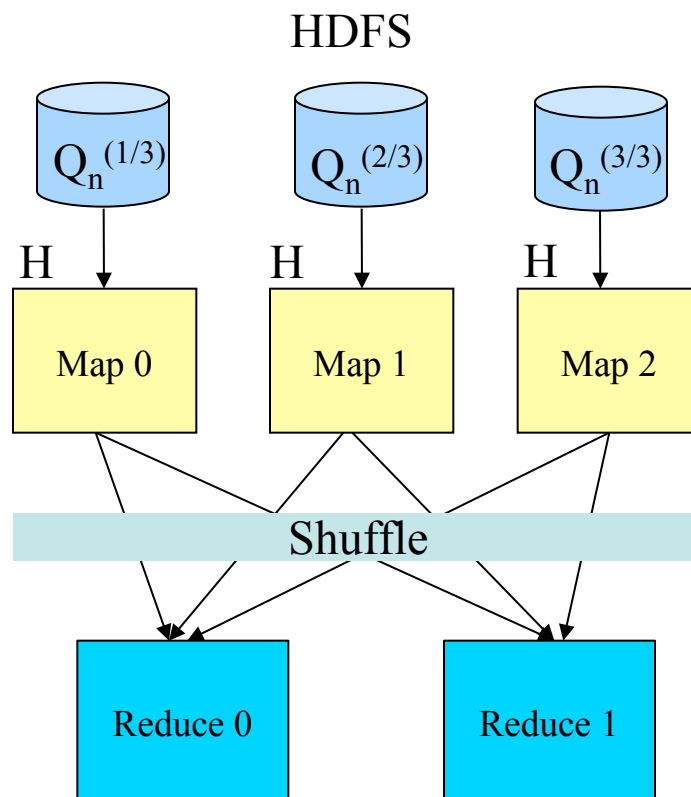
Q_n : O(100 Gbytes)
 H : O(Kbytes)

- Naïve multiplication: too expensive
- Proposed:
 - 'cache'-based multiplication: broadcast the small matrix H to all the machines that contains Q_n

Skewed Matrix-Matrix Mult.

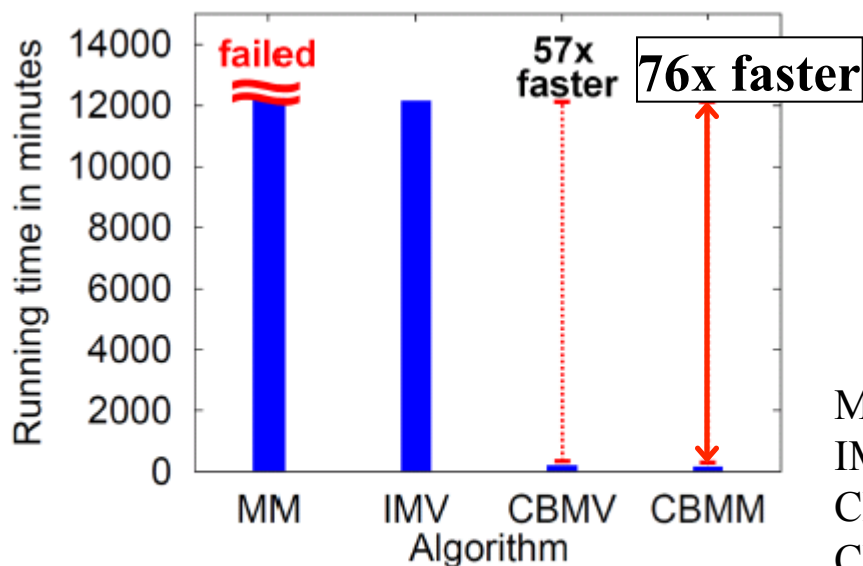
Details

- **'cache'-based** multiplication: broadcast the small matrix H to all the machines that contains Q_n

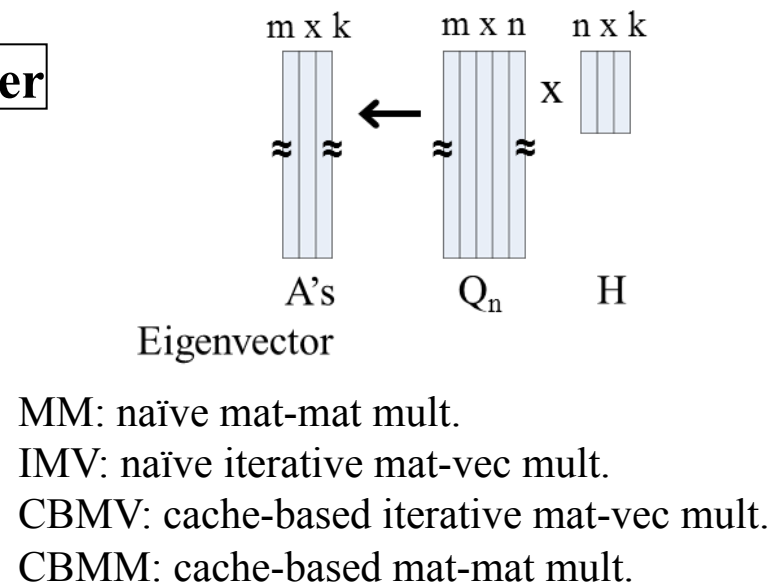


Skewed Matrix-Matrix Mult.

Which Matrix-Matrix multiplication algorithm runs the fastest?



Time vs. algorithms



(100 machines used)

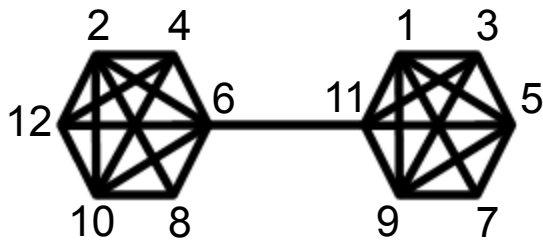
Cache-based MM runs 76x faster

Outline

- Problem#1: Patterns in graphs
- Problem#2: Tools
- Problem#3: Scalability - PEGASUS
 - Structure Analysis
 - Eigensolver
 - ➔ – Graph Layout and Compression
- Conclusions

Node Order Matters

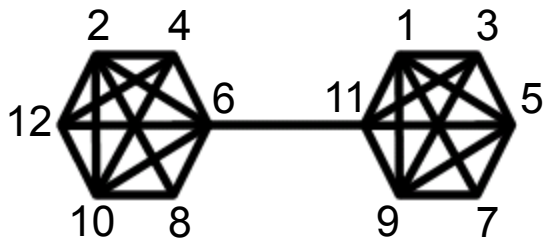
- A graph and the adjacency matrix



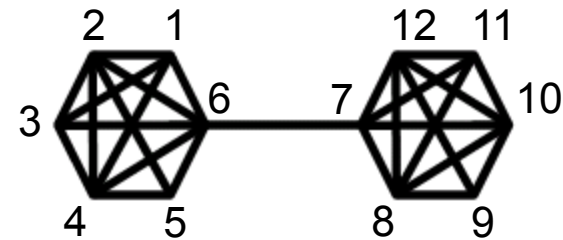
	1	2	3	4	5	6	7	8	9	10	11	12
1			1	1	1	1	1	1	1	1		
2			1	1	1	1	1	1	1	1		
3	1	1			1	1	1	1	1	1		
4	1	1			1	1	1	1	1	1		
5	1	1	1	1			1	1	1	1		
6	1	1	1	1			1	1	1	1	1	
7	1	1	1	1	1	1			1	1	1	
8	1	1	1	1	1	1			1	1	1	
9	1	1	1	1	1	1	1	1			1	
10	1	1	1	1	1	1	1	1			1	
11	1	1	1	1	1	1	1	1	1	1		
12	1	1	1	1	1	1	1	1	1	1		

Node Order Matters

- Same graphs with different orderings



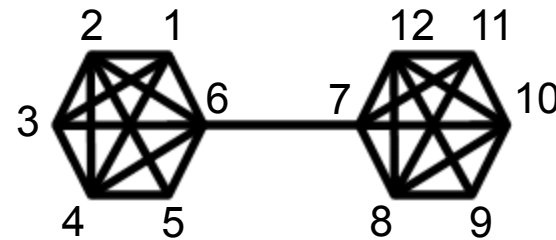
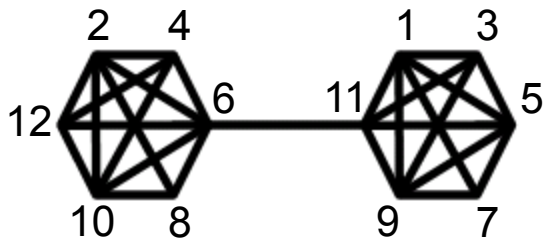
	1	2	3	4	5	6	7	8	9	10	11	12
1			1	1	1	1	1	1	1	1		
2			1	1	1	1	1	1	1	1		
3	1	1			1	1	1	1	1	1		
4	1	1			1	1	1	1	1	1		
5	1	1	1			1	1	1	1	1		
6	1	1	1			1	1	1	1	1		
7	1	1	1	1				1	1	1		
8	1	1	1	1				1	1	1		
9	1	1	1	1	1	1			1	1		
10	1	1	1	1	1	1			1	1		
11	1	1	1	1	1	1	1					
12	1	1	1	1	1	1	1					



	1	2	3	4	5	6	7	8	9	10	11	12
1		1	1	1	1	1						
2	1		1	1	1	1						
3	1	1		1	1	1						
4	1	1	1		1	1						
5	1	1	1	1		1						
6	1	1	1	1	1		1					
7						1		1	1	1	1	1
8							1		1	1	1	1
9								1		1	1	1
10									1	1	1	1
11										1	1	1
12											1	1

Good ordering = Good compression

- Same graphs with different orderings



Many sparse blocks

	1	2	3	4	5	6	7	8	9	10	11	12
1			1		1		1		1		1	
2				1		1		1		1		1
3	1				1		1		1		1	
4		1				1		1		1		1
5	1		1				1		1		1	
6		1		1				1		1	1	1
7	1			1					1		1	
8		1			1					1		1
9	1		1			1		1			1	
10		1		1			1					1
11	1		1		1	1	1		1			
12		1		1		1		1				



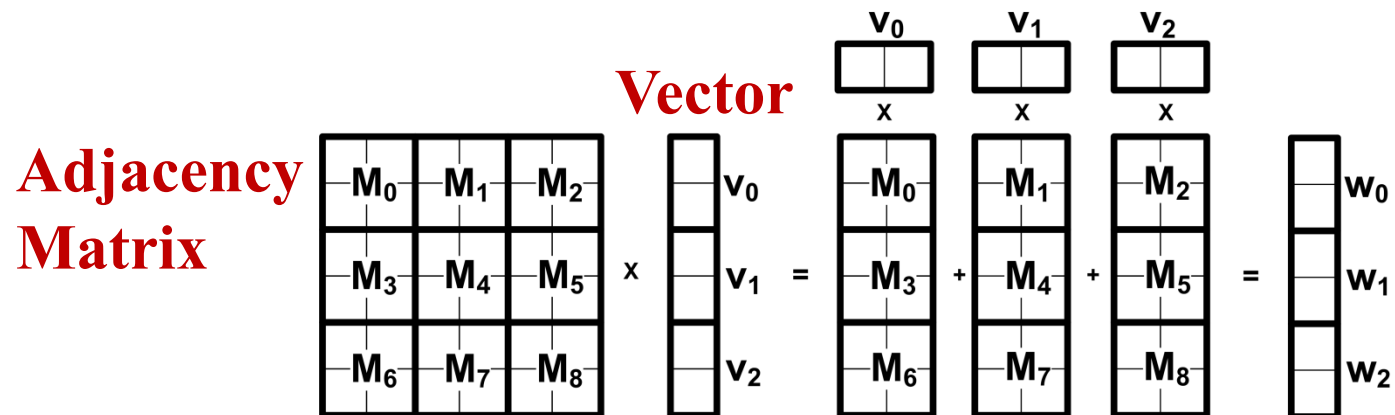
Few dense blocks

	1	2	3	4	5	6	7	8	9	10	11	12
1		1	1	1	1	1						
2	1		1	1	1	1						
3	1	1		1	1	1						
4	1	1	1		1	1						
5	1	1	1	1		1						
6	1	1	1	1	1		1					
7						1		1	1	1	1	1
8							1		1	1	1	1
9								1		1	1	1
10									1	1	1	1
11										1	1	1
12											1	1



Application

- Block-based matrix-vector multiplication



Few, dense blocks

\Rightarrow Better compression, faster running time

Problem Definition

- Given a graph, how can we lay-out its edges so that nonzero elements are well-clustered?
- Better clustering = better compression

Many
sparse
blocks



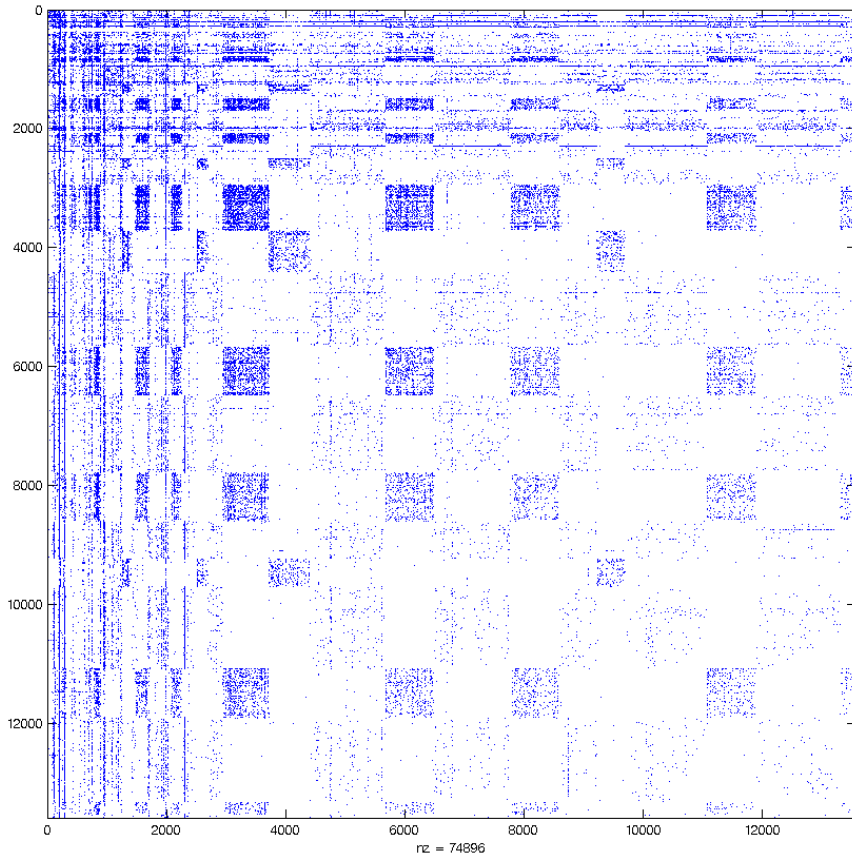
	1	2	3	4	5	6	7	8	9	10	11	12
1			1		1		1		1		1	
2				1		1		1		1		1
3	1				1		1		1		1	
4		1				1		1		1		1
5	1		1				1		1		1	
6		1		1				1		1	1	1
7	1			1		1			1		1	
8		1			1				1		1	
9	1		1		1		1				1	
10		1		1		1		1				1
11	1		1		1	1	1		1			
12		1		1		1		1		1		

Few
dense
blocks

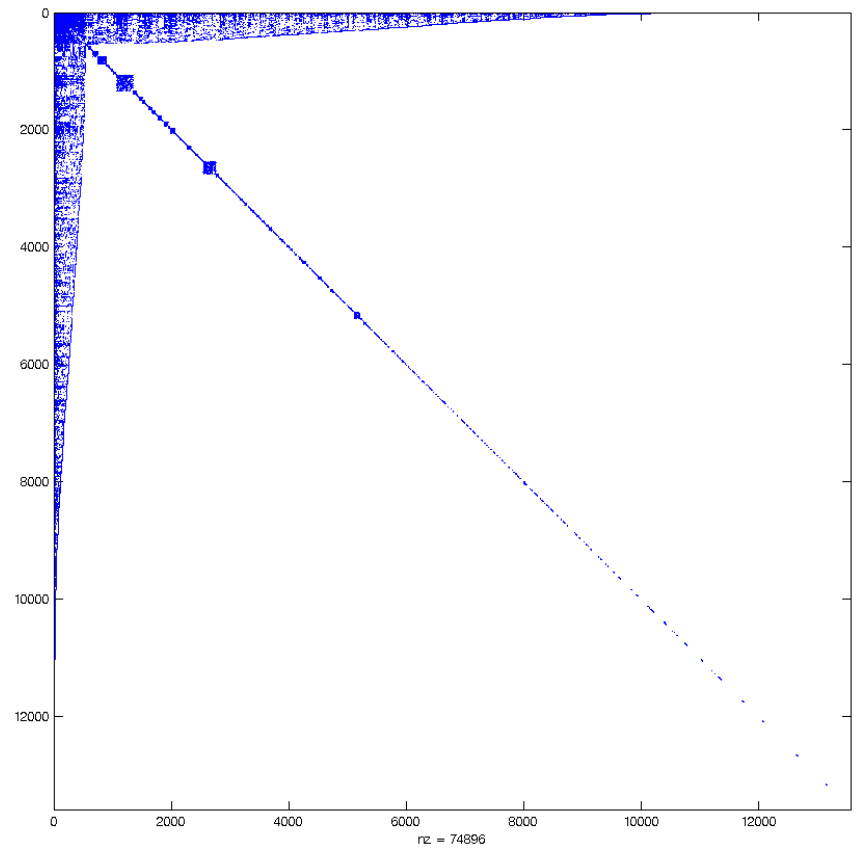


	1	2	3	4	5	6	7	8	9	10	11	12
1		1	1	1	1	1						
2	1		1	1	1	1						
3	1	1		1	1	1						
4	1	1	1		1	1						
5	1	1	1	1		1						
6	1	1	1	1	1		1					
7						1		1	1	1	1	1
8							1		1	1	1	1
9								1		1	1	1
10									1	1	1	1
11										1	1	1
12											1	1

Main Result



Original



SlashBurn

Outline

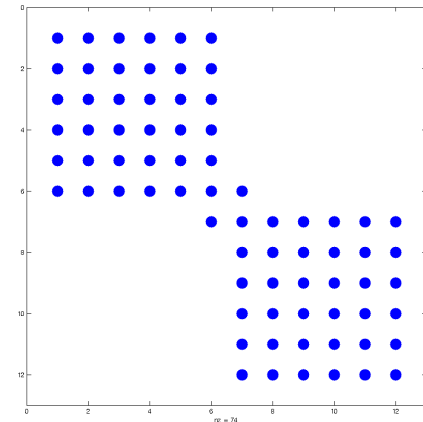
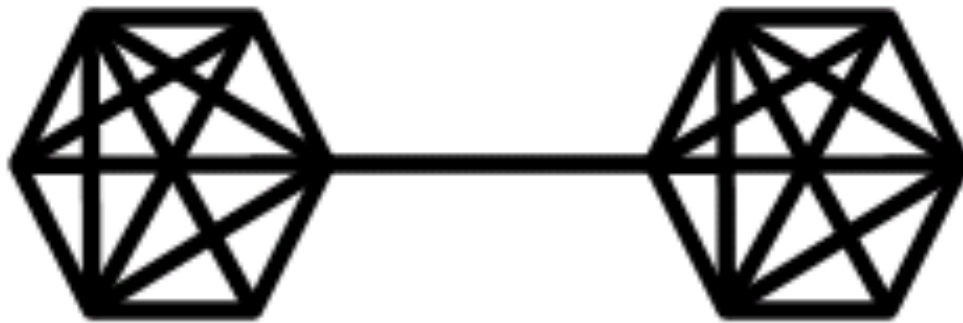
- ...
- Problem#3: Scalability - PEGASUS
 - Structure Analysis
 - Eigensolver
 - Graph Layout and Compression
 - ➔ • Proposed Method
 - Results
- Conclusions

Survey

- Given a graph, how can we lay-out its edges so that nonzero elements are well-clustered?
 - 1) Graph based clustering
 - Normalized cut, spectral clustering
 - 2) Heuristics
 - Lexicographic ordering for Web
 - Shingle ordering

1) Graph Based Clustering

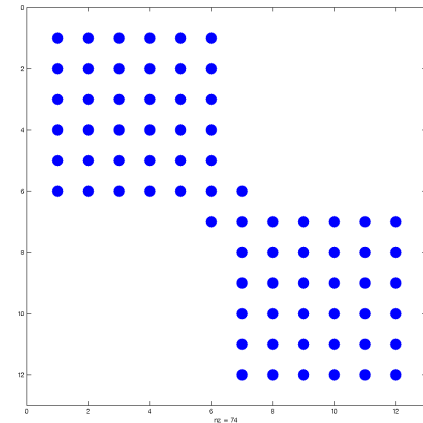
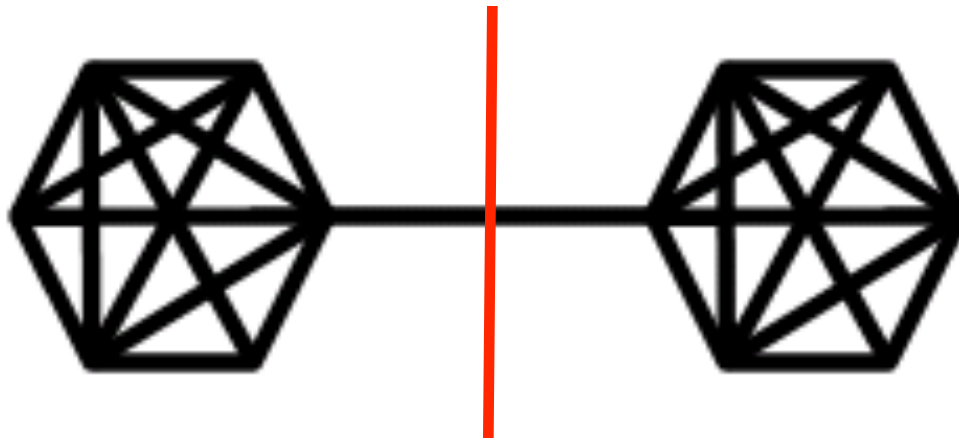
- Goal: find homogeneous sets of nodes from graph S
 - E.g.) Spectral clustering and normalized cut
 - Many intra-edges, few inter-edges



Caveman Communities

1) Graph Based Clustering

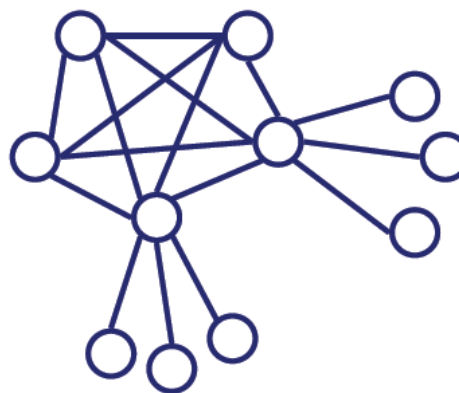
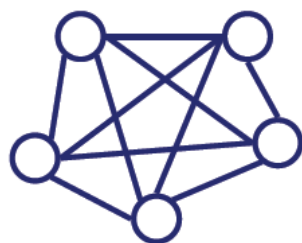
- Goal: find homogeneous sets of nodes from graph
 - E.g.) Spectral clustering and normalized cut
 - Many intra-edges, few inter-edges



Caveman Communities

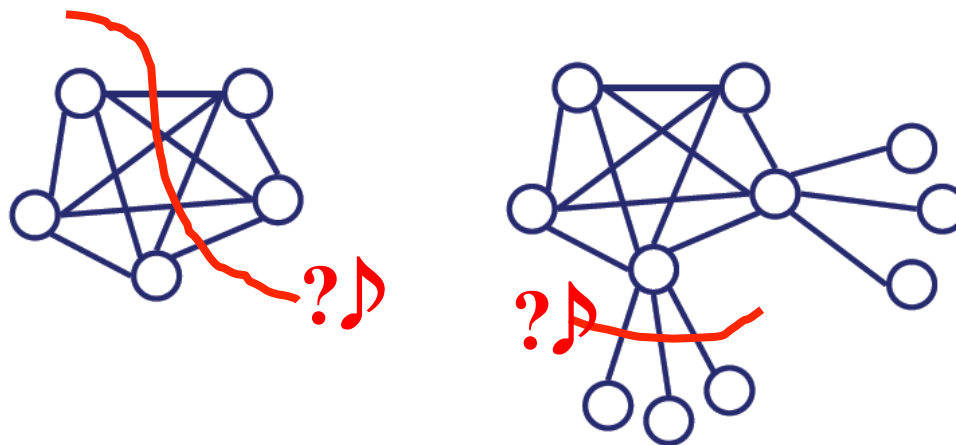
1) Graph Based Clustering

- But, real graphs: no good cuts
 - [Tauro+ 01], [Siganos+ 06], [Chakrabarti +04], [Leskovec+ 08]



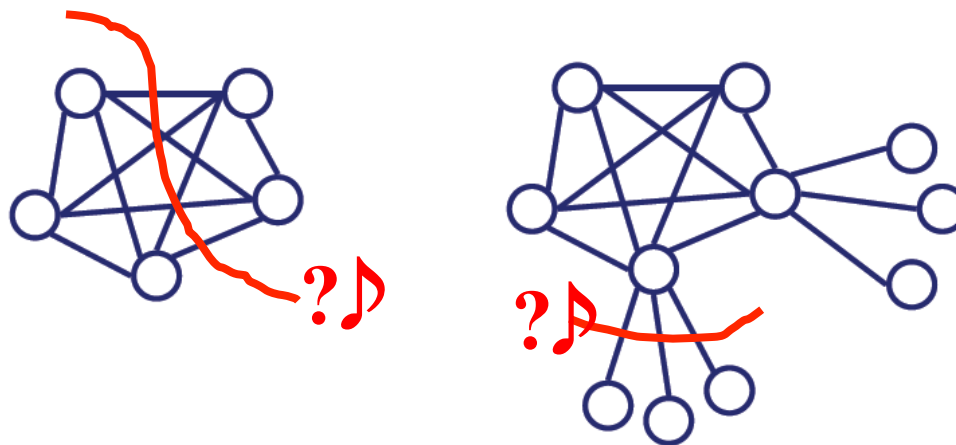
1) Graph Based Clustering

- But, real graphs: no good cuts
 - [Tauro+ 01], [Siganos+ 06], [Chakrabarti +04], [Lesko vec+ 08]



1) Graph Based Clustering

- But, real graphs: no good cuts
 - [Tauro+ 01], [Siganos+ 06], [Chakrabarti +04], [Lesko vec+ 08]



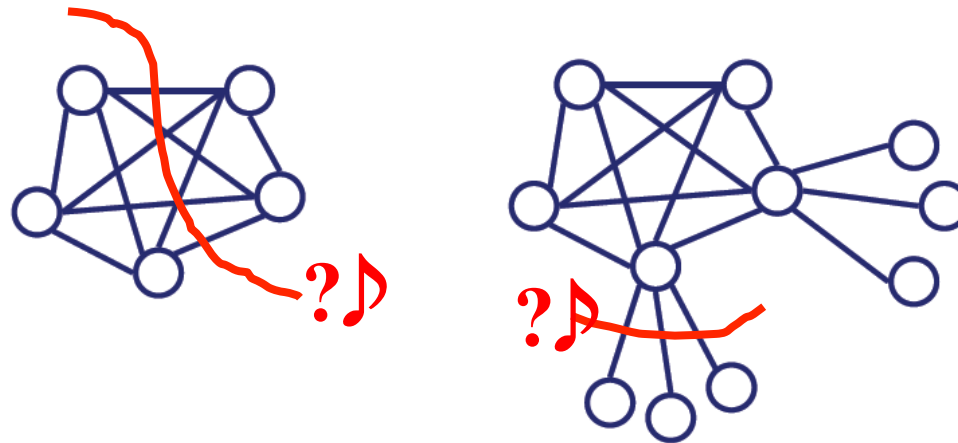
- What should we do?

2) Older Heuristics

- Web graph: lexicographic ordering [Boldi+, 04]
 - Locality : many intra edges between neighbors
 - Similarity : out links of neighbors are similar
- Social network: shingle ordering [Chierichetti+ 09]
 - Group nodes with similar out-neighbors

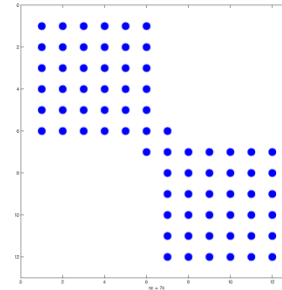
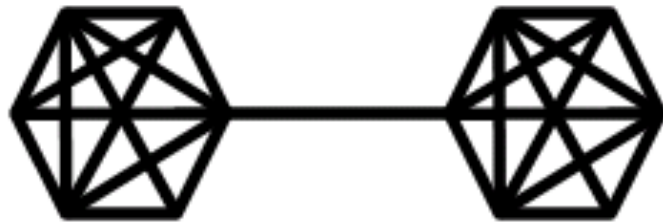
Summary : Previous Work

- Tries to find homogenous regions for graph compression
 - Fails to find them, because they often don't exist



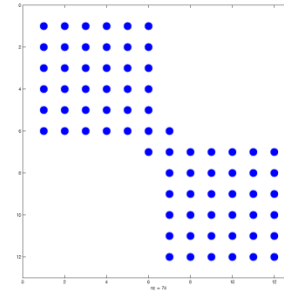
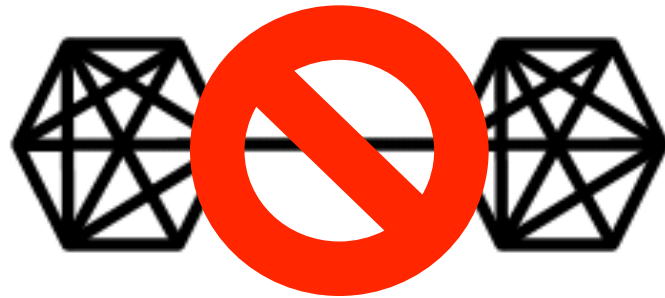
Our Observation

- Caveman assumption



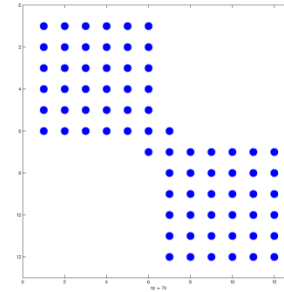
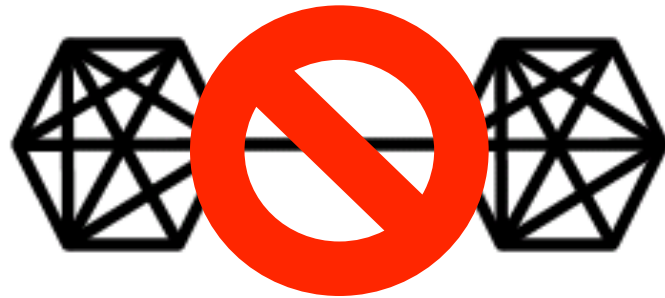
Our Observation

- Caveman assumption: wrong!

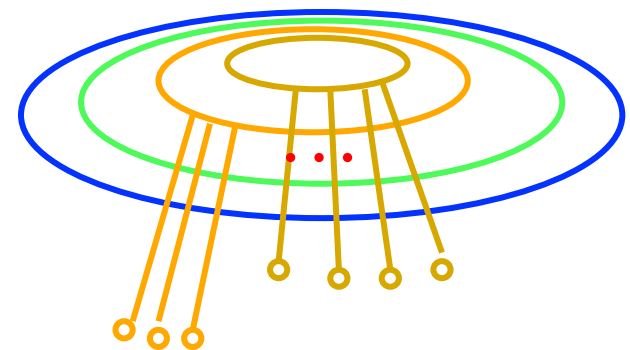


Our Solution

- Caveman assumption: wrong!

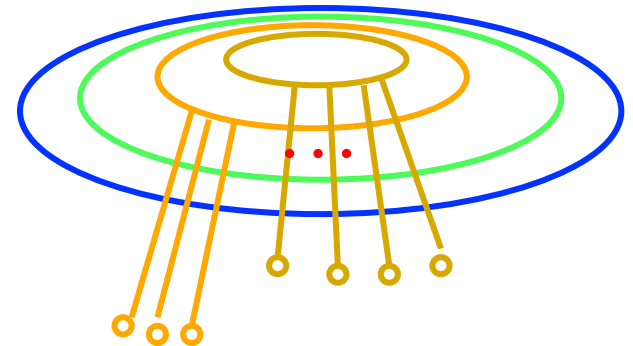


- Instead, we envision graphs as **nodes** connected by **connectors** connected by **super connectors**...



Our Solution

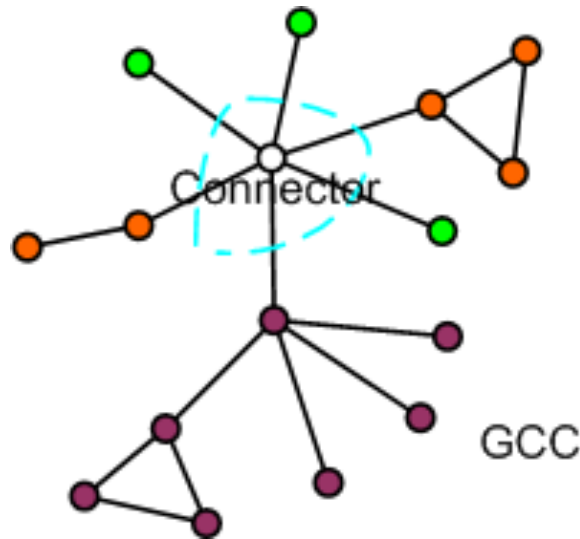
- Instead, we envision graphs as **nodes** connected by **connectors** connected by **super connectors**...
- Use “Graph Shattering” to ‘peel’ the graphs from **super connectors**, and then **connectors**,
...♪



Graph Shattering

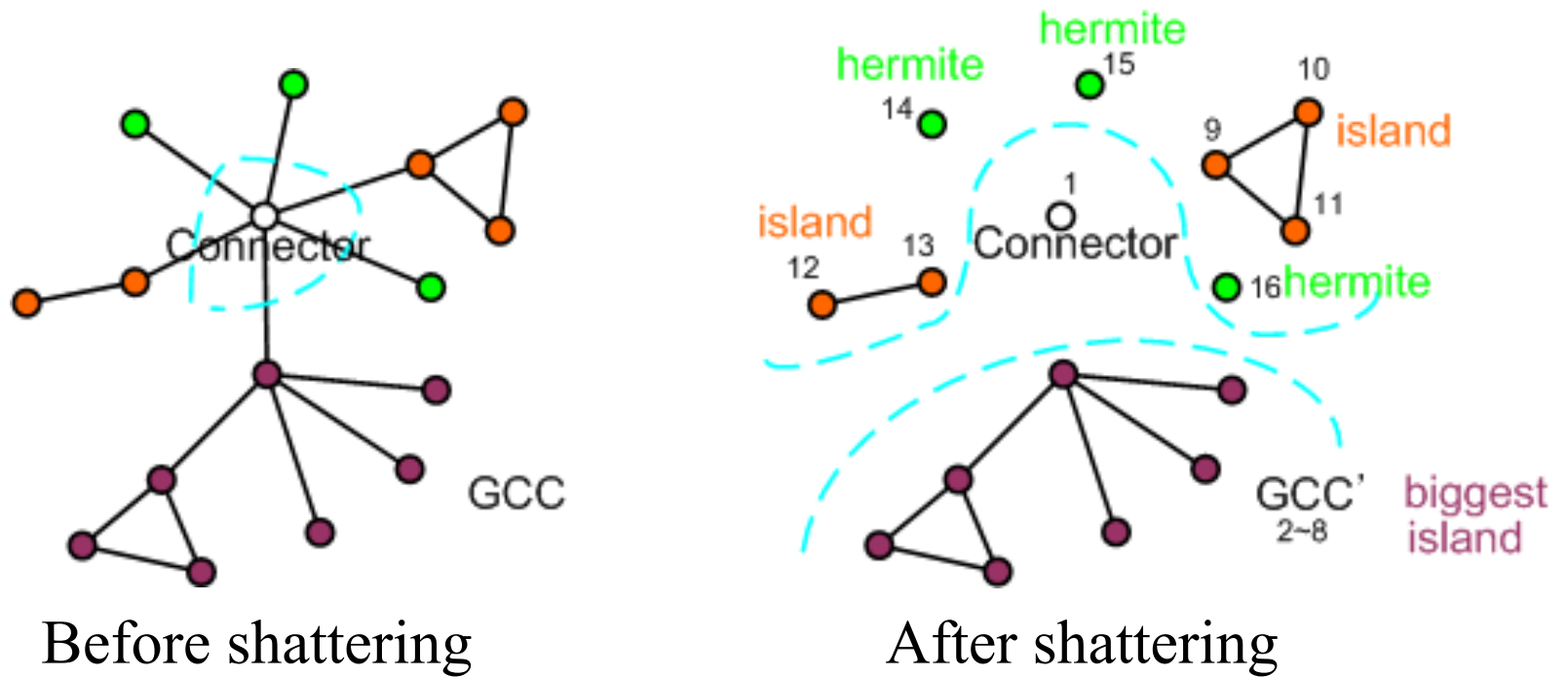
- k -shattering of a graph G
 - Removes top k connectors and their incident edges from G

Graph Shattering

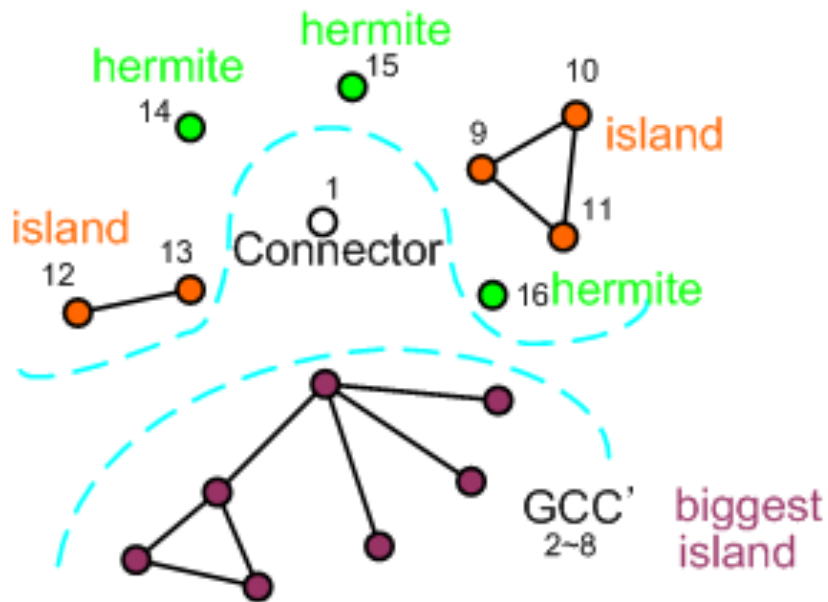


Before shattering

Graph Shattering



Graph Shattering

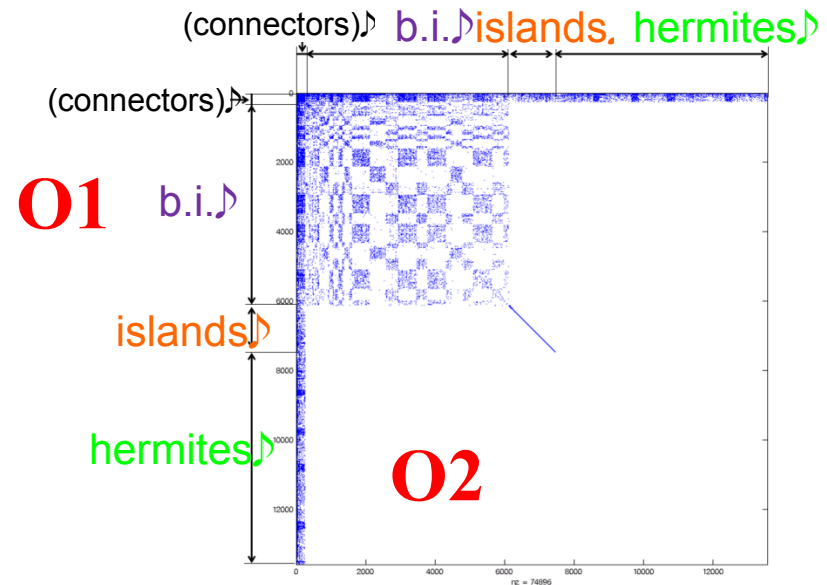


After shattering

■ Observations in real graphs

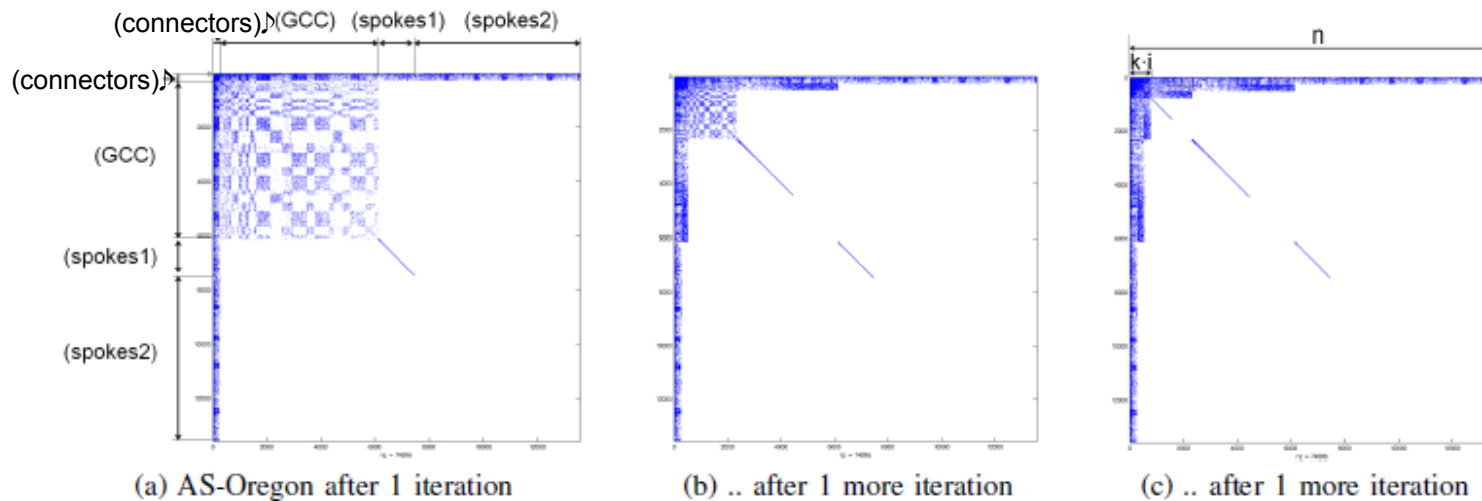
O1. Portion of GCC is much smaller after shattering

O2. A lot of disconnected components



Slash-Burn method (intuition)

- ‘burn’ the top k connectors, and ‘slash’ the edges
- Move k connectors to the front of the row/column, sort connected components by decr. size
- Recurse on the remaining GCC



Outline

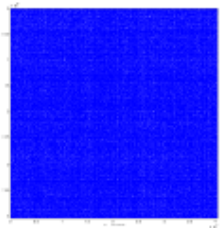
- ...
 - Problem#3: Scalability - PEGASUS
 - Structure Analysis
 - Eigensolver
 - Graph Layout and Compression
 - Proposed Method
 - Results
- ➔
- Conclusions

Goal of Experiments

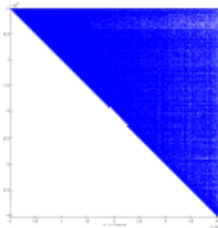
- [Q1] Compression savings?
- [Q2] Running time savings?

A1. Compression

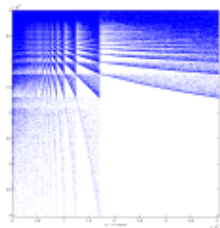
Winner



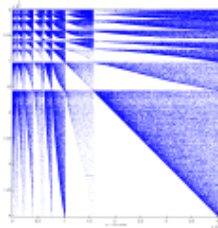
Flickr:
(a) Random



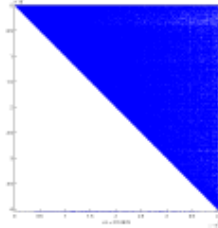
(b) Natural



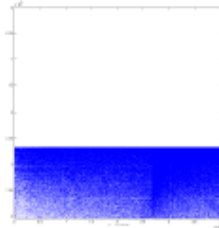
(c) Degree Sort



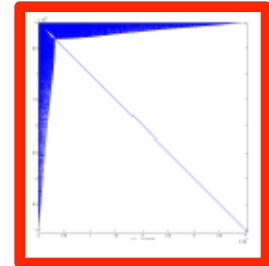
(d) Cross Association



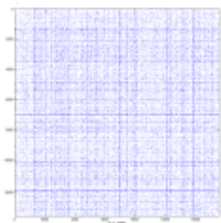
(e) Spectral Clustering



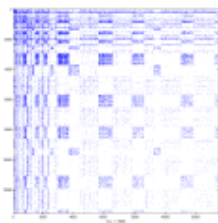
(f) Shingle



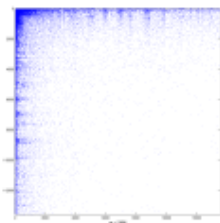
(g) SLASHBURN



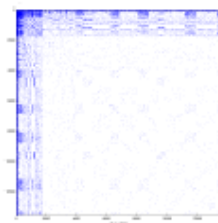
AS-Oregon:
(a) Random



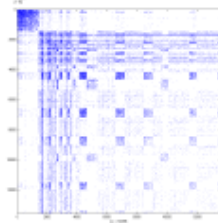
(b) Natural



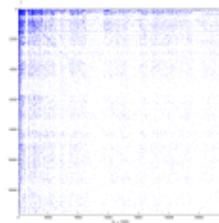
(c) Degree Sort



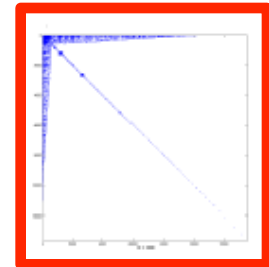
(d) Cross Association



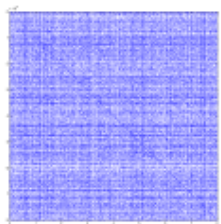
(e) Spectral Clustering



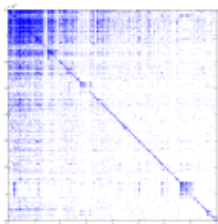
(f) Shingle



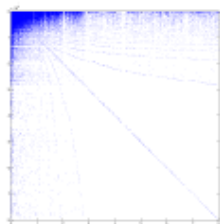
(g) SLASHBURN



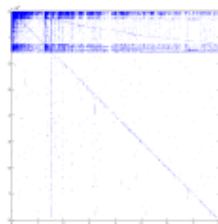
Enron:
(a) Random



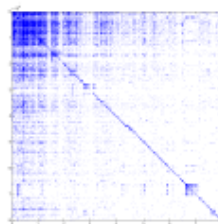
(b) Natural



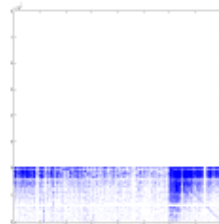
(c) Degree Sort



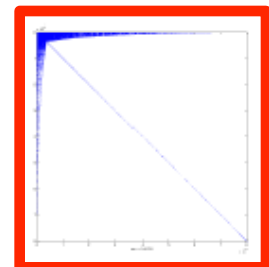
(d) Cross Association



(e) Spectral Clustering



(f) Shingle



(g) SLASHBURN

A1. Compression

Details

■ Cost functions used

- 1) Number of non-empty blocks
- 2) Information theoretic cost : minimum bits to encode non-zero elements inside blocks

model complexity

$$|T| \cdot 2 \log \frac{n}{b} + \sum_{\tau \in T} b^2 \cdot H\left(\frac{z(\tau)}{b^2}\right)$$

costs given the model

$|T|$: # of nonempty blocks

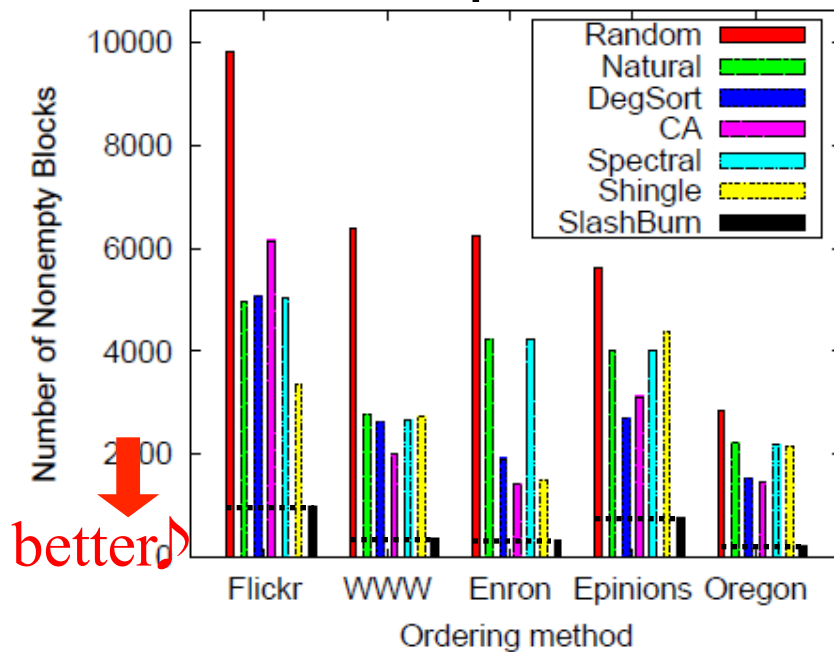
n : # of nodes

b : block width

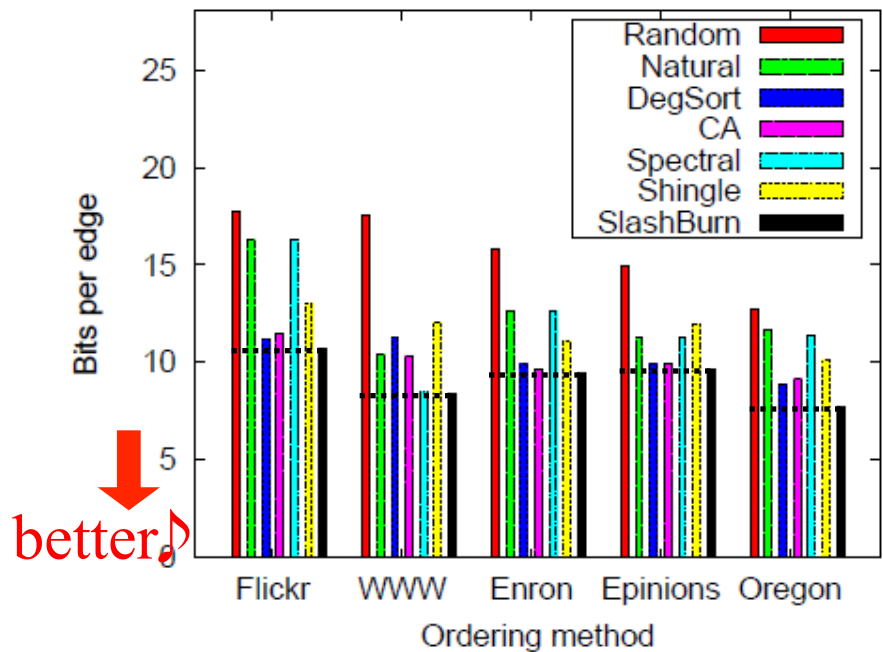
$H()$: Shannon entropy func

A1. Compression

■ Cost comparison



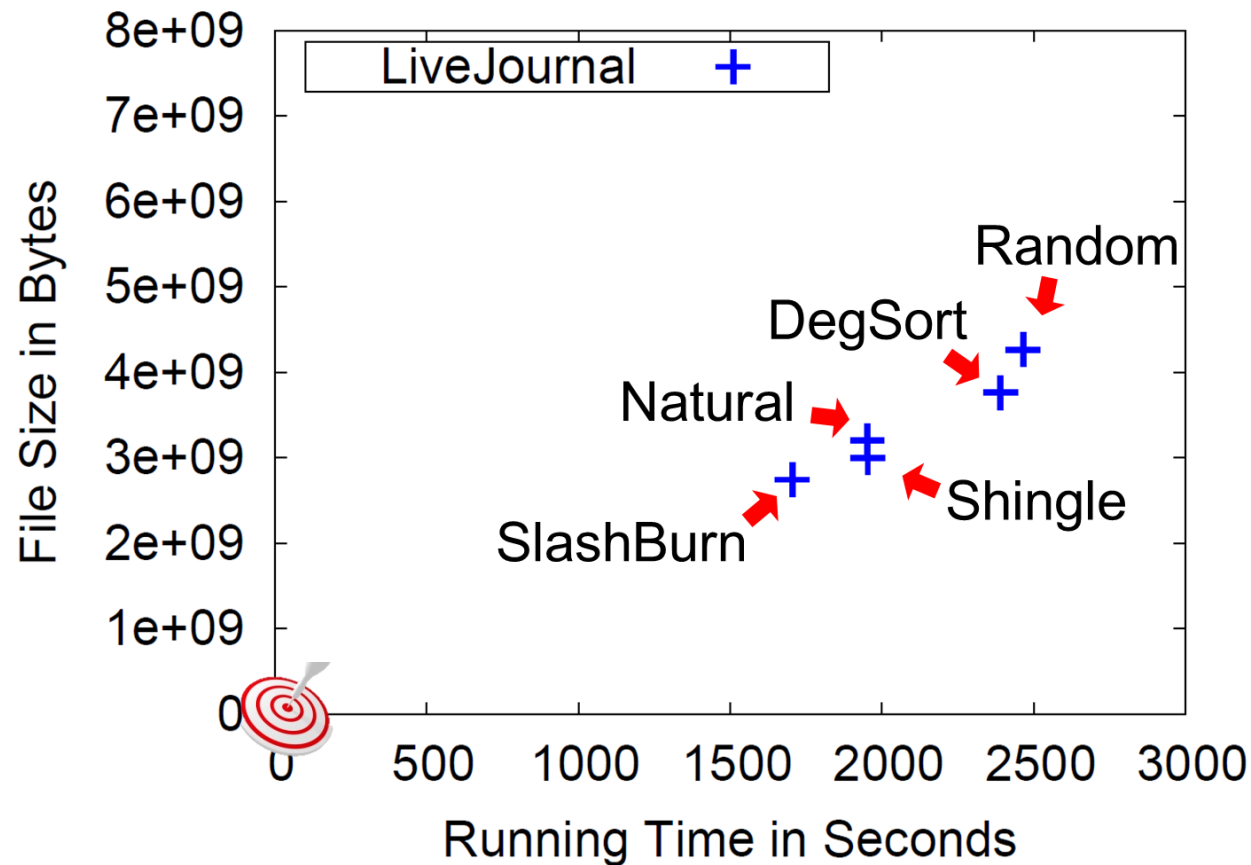
(a) $cost_{nz}(A, b)$: number of nonempty blocks



(b) $cost_{it}(A, b)$: information theoretic cost

- SlashBurn outperforms all competitors for all dataset!
(smallest number of nonempty blocks, as well as bits per edge)

A2. Running Time



- SlashBurn outperforms all competitors !
(running time as well as file size)

Outline

- Introduction – Motivation
- Problem#1: Patterns in graphs
- Problem#2: Tools
- Problem#3: Scalability - PEGASUS
- ➔ • Conclusions

Conclusions

- PEGASUS: Peta-Scale Graph Mining System
 - Patterns and anomalies in large graphs
 - PageRank, Connected Components, Radius, Eigensolver
 - Outreach
 - Microsoft : part of Hadoop distribution for Windows Azure
 - One of the core systems for several DARPA projects (ADA MS, INARC, DTRA)



www.cs.cmu.edu/~pegasus



Conclusions

- High impact applications require large graph mining

Cyber
Security



Fraud
Detection



Social
Network

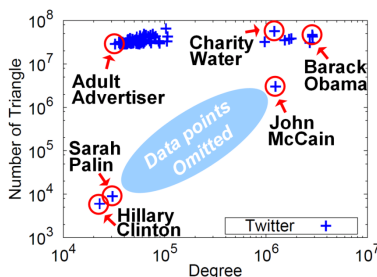
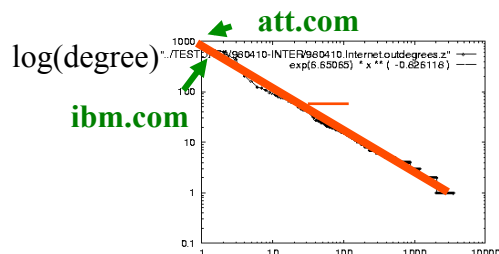


Search
Engine



Health
Care





www.cs.cmu.edu/~pegasus

Thank you !

Complementary tutorial: *Mining Billion-Scale Graphs: Systems and Implementations*: Haixun Wang et al